



GLOBAL JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY: C  
SOFTWARE & DATA ENGINEERING  
Volume 16 Issue 5 Version 1.0 Year 2016  
Type: Double Blind Peer Reviewed International Research Journal  
Publisher: Global Journals Inc. (USA)  
Online ISSN: 0975-4172 & Print ISSN: 0975-4350

## Security in Data Mining- A Comprehensive Survey

By Niranjana A, Nitish A, P Deepa Shenoy & Venugopal K R

*University Visvesvaraya College of Engineering*

**Abstract-** Data mining techniques, while allowing the individuals to extract hidden knowledge on one hand, introduce a number of privacy threats on the other hand. In this paper, we study some of these issues along with a detailed discussion on the applications of various data mining techniques for providing security. An efficient classification technique when used properly, would allow an user to differentiate between a phishing website and a normal website, to classify the users as normal users and criminals based on their activities on Social networks (Crime Profiling) and to prevent users from executing malicious codes by labelling them as malicious. The most important applications of Data mining is the detection of intrusions, where different Data mining techniques can be applied to effectively detect an intrusion and report in real time so that necessary actions are taken to thwart the attempts of the intruder.

**Keywords :** *anomaly detection; classification; intrusion detection system; outlier detection; privacy preserving data mining.*

**GJCST-C Classification :** *H.2.8*



*Strictly as per the compliance and regulations of:*



# Security in Data Mining- A Comprehensive Survey

Niranjan A <sup>α</sup>, Nitish A <sup>σ</sup>, P Deepa Shenoy <sup>ρ</sup> & Venugopal K R <sup>ω</sup>

**Abstract** Data mining techniques, while allowing the individuals to extract hidden knowledge on one hand, introduce a number of privacy threats on the other hand. In this paper, we study some of these issues along with a detailed discussion on the applications of various data mining techniques for providing security. An efficient classification technique when used properly, would allow an user to differentiate between a phishing website and a normal website, to classify the users as normal users and criminals based on their activities on Social networks (Crime Profiling) and to prevent users from executing malicious codes by labelling them as malicious. The most important applications of Data mining is the detection of intrusions, where different Data mining techniques can be applied to effectively detect an intrusion and report in real time so that necessary actions are taken to thwart the attempts of the intruder. Privacy Preservation, Outlier Detection, Anomaly Detection and PhishingWebsite Classification are discussed in this paper.

**Keywords:** anomaly detection; classification; intrusion detection system; outlier detection; privacy preserving data mining.

## 1. INTRODUCTION

The term Security from the context of computers is the ability, a system must possess to protect data or information and its resources with respect to confidentiality, integrity and authenticity[1]. Confidentiality ensures that, a third party in no way would be able to read and understand the content while Integrity would not allow a third party to change or modify the content as a whole or even parts of it. Authenticity feature on the other hand would not allow a person to use, view or modify the content or the resource, if he is found to be unauthorised[2].

Those actions that compromise the availability, integrity or confidentiality of one or more resources of a computer could be termed as *Intrusion*. Preventing intrusions employing firewall and filtering router policies fail to stop these attacks. In spite of all attempts to build secure systems, intrusions can still happen and hence they must be detected on their onset. An Intrusion detection system(IDS)[3] by employing data mining techniques can discover consistent patterns of features of a system that are useful can detect anomalies and known intrusions using a relevant set of classifiers. Using some of the basic data mining techniques such

as Classification and Clustering, Intrusion can be detected easily. Classification techniques are helpful in analyzing and labelling the test data into known type of classes, while Clustering techniques are used to group objects into a set of clusters, such that all similar objects become the members of the same cluster and all other objects become members of other clusters[4]. Data mining, while allowing the extraction of hidden patterns or the underlying

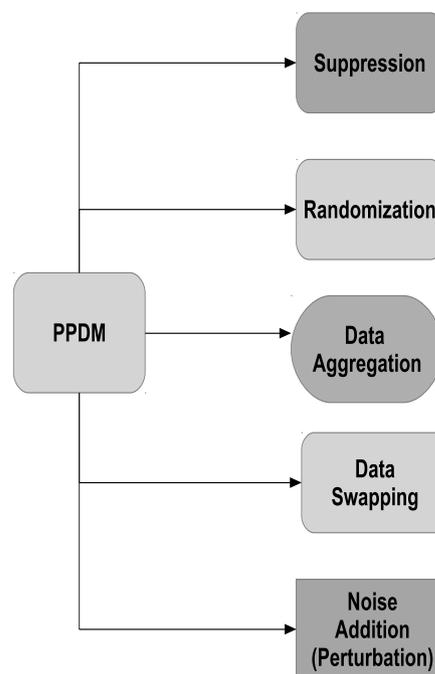


Figure 1: Privacy Preserving Data Mining Techniques

knowledge from large volumes of data, might pose security challenges[5]. Privacy Preserving Data Mining(PPDM)aims at safeguarding sensitive information from an un-solicited or unsanctioned disclosure[6]. A number of PPDM approaches have been proposed so far. Some of them are listed as shown in Fig. 1, based on their enforcing privacy principle.

### a) Suppression

Any private or sensitive information pertaining to an individual such as name, age, salary, address and other information is suppressed before any computation takes place. Some of the techniques employed for this suppression are Rounding(Rs/- 35462.33 may be rounded to 35,000), Generalization (Name Louis Philip

Author <sup>α</sup> <sup>σ</sup> <sup>ρ</sup> <sup>ω</sup>: Department of Computer Science and Engineering, University Visvesvaraya College of Engineering, Bangalore University, Bangalore, India. e-mail: niranjan.a.in@ieee.org

may be replaced with the initials LP and Place Hamburg may be replaced with HMG and so forth). However when data mining requires full access to sensitive values, Suppression cannot be used. An alternate way of suppression is to limit the identity linkage of a record rather than suppressing the sensitive information present within a record. This technique is referred to as *De-Identification*. *k-Anonymity* is one such de-identification technique. It ensures that protection of the data released against *Re-identification* of the persons to which the data refer [7][8]. Enforcing *k-anonymity* before all data are collected in one trusted place is difficult. A cryptographic solution based on Secret Sharing technique of Shamir could be used instead; this however incurs computation overhead.

#### b) *Randomization*

Assuming the presence of a central server of a company that accepts information present with many customers and performs data mining techniques for building an Aggregate Model; *Randomization* allows the customers to introduce controlled noise or randomly perturb the records and to take away true information present in it. Introduction of noise can be achieved in several ways by addition or multiplication of the values generated randomly. *Perturbation* helps *Randomization* technique to achieve preservation of the required privacy.

The individual records are generated by the addition of such randomly generated noise to the original data. The noise thus added to individual records cannot be recovered, resulting in the desired privacy. *Randomization* techniques typically involve the following steps:

1. Only after randomizing their data, the Data Providers transmit this data to the Data Receiver.
2. Data receiver computes the distribution by running a Distribution Reconstruction Algorithm.

#### c) *Data Aggregation*

*Data Aggregation* Techniques, in order to facilitate data analysis: combine data together from various sources. This might allow an attacker to deduce private and individual-level data and to identify the party. When the extracted data allows the data miner to identify specific individuals, his privacy is considered to be under a serious threat. To prevent data from being identified, it may be anonymized immediately after the aggregation process. However, the Anonymized data sets can still contain enough information that could be used for the identification of individuals [9].

#### d) *Data Swapping*

*Data swapping* process involves swapping of values across different records for the sake of privacy-preservation. Without perturbing the lower order totals of the data, privacy of data can still be preserved allowing aggregate computations to be performed exactly as before. Since this technique does not follow randomization, it can be used in conjunction with other

frameworks such as *k-anonymity* without violating the privacy definitions for that model.

#### e) *Noise Addition/Perturbation*

Differential privacy through the addition of controlled noise provides a mechanism that maximizes the accuracy of queries while minimizing the chances of identification of its records [10]. Some of the techniques used in this regard are:

1. Laplace Mechanism
2. Sequential Composition
3. Parallel Composition

The rest of this paper is structured as follows: Section- II covers a brief review of Classification and Detection of intrusions by employing various Data Mining Techniques, while Clustering techniques and their applications in Intrusion Detection are presented in Section-III. PPDM techniques and their necessity along with various types of PPDM are discussed in Section-IV. An overview of Intrusion Detection System is discussed in Section-V. Phishing Website Classification using Data Mining Techniques are presented in Section-VI. Artificial Neural Networks(ANN) are presented in Section-VII. Section- VIII presents Anomaly Detection/Outlier Detection. Section- IX describes the various ways of Mitigating Code Injection Attacks.

## II. CLASSIFICATION AND DETECTION USING DATA MINING TECHNIQUES

Malware computer programs that replicate themselves in order to spread from one computer to another computer are called as worms. Malware includes worms, computer viruses, Trojan Horse, key loggers, adware, spyware Port scan worm, UDP worm, http worm, User to Root Worm and Remote to Local Worm and other malicious code [11]. Attackers write these programs for various reasons varying from interruption of a computer process, gathering sensitive information, or gaining entry to private systems. Detecting a worm on the internet is very important, because it creates vulnerable points and reduces the performance of the system. Hence it is essential to detect the worm on the onset and classify it using data mining classification algorithms much before it causes any damage.

Some of the classification algorithms that can be used are Random Forest, Decision Tree, Bayesian and others [12]. A majority of worm detection techniques use Intrusion Detection System (IDS) as the underlying principle. Automatic detection is challenging because it is tough to predict what form the next worm will take. IDS can be classified into two types namely Network based IDS and Host based IDS. The Network based Intrusion Detection System reflects network packets before they spread to an end-host, while the Host based Intrusion Detection System reflects network packets that are already spread to the end-host. Moreover, the Host based detection studies encode network packets so

that the stroke of the internet worm may be struck. When we focus on the network packet without encoding, we must study the performances of traffic in the network. Several machine learning techniques have been used in the field of intrusion and worm detection systems. Thus, Data Mining and in particular Machine Learning Technique has an important role and is essential in worm detection systems. Using various Data Mining schemes several new techniques to build several Intrusion Detection models have been proposed. Decision Trees and Genetic Algorithms of Machine Learning can be employed to learn anomalous and normal patterns from the training set and classifiers are then generated based on the test data to label them as Normal or Abnormal classes. The data that is labelled as Abnormal could be a pointer to the presence of an intrusion.

a) *Decision Trees*

Quinlan's decision tree technique, is one of most popular machine learning techniques. The tree is constructed using a number of decision and leaf nodes following divide-and-conquer technique[12]. Each decision node tests a condition on one of the attributes of the input data and can essentially have a number of branches, to handle a separate outcome of the test. The result of decision may be represented as a leaf node. A training data set  $T$  is a set of  $n$ -classes  $\{C1, C2, \dots, Cn\}$ .  $T$  is treated as a leaf when it comprises of cases belonging to a single class. If  $T$  is empty with no cases, it is still treated a leaf and the major class of the parent node is given the related class. A test based on an attribute  $a_i$  of the training data is performed when  $T$  consists of multiple classes,  $T$  is split into  $k$  subsets  $\{T1, T2, \dots, Tk\}$ , where  $k$  gives the number of test outcomes. The process is recursed over each  $T_j$ , where  $1 \leq j \leq n$ , until every subset belongs to a single class. Choosing the best attribute for each decision node while constructing the decision tree is very crucial. The C4.5-DT adopts Gain Ratio Criterion for the same. According to this criterion, an attribute that provides maximum information gain and that reduces the bias in favor of tests is chosen. The tree thus built can then be used to classify the test data, whose features are same as that of the training data. The test is carried out starting from the root node. Based on the outcome, one of the branches leading to a child is followed. As long as the child is not a leaf, the process is repeated recursively. The class and its corresponding leaf node is given to the test case being examined.

b) *Genetic Algorithms(GA)*

A machine learning approach of solving problems by employing biological evolution techniques are called *Genetic Algorithms(GA)*. They can be effectively used to optimize a population of candidate solutions. GA makes use of data structures that are modelled on chromosomes and they are subjected to

evolution using genetic operators namely: selection, crossover and mutation[13]. Random generation of a population of chromosomes is performed in the beginning. The population thus formed comprises of all possible solutions of a problem and are considered the candidate solutions. Different positions of a chromosome called '*genes*' are encoded as bits, characters or numbers. A function called *Fitness Function* evaluates the goodness of each chromosome based on the desired solution. *Crossover* operator simulates natural reproduction while *Mutation* operator simulates mutation of the species. The *Selection* operator chooses the fittest chromosomes[14]. Fig 2. depicts the operations of Genetic Algorithms. Before using GA for solving various problems, following three factors have to be considered

1. Fitness function
2. Individuals representation and
3. Parameters of GA

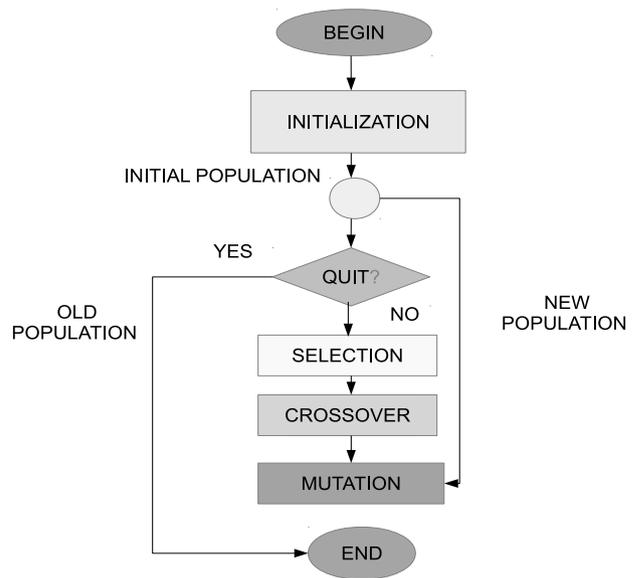


Figure 2: Flowchart for a GA

GA based approach can be incorporated for designing Artificial Immune Systems. Using this approach, Bin et al.,[15] have proposed a method for smartphone malware detection where static and dynamic signatures of malwares are extracted and malicious scores of tested samples are obtained.

c) *Random Forest*

A classification algorithm that is made up of a collection of tree structured classifiers, and that chooses the winner class based on the votes casted by the individual trees present in the forest is called the *Random Forest Algorithm*. Each tree is constructed by picking up random data from a training dataset. The chosen dataset may be split up into training and testsets. The major chunk of the dataset goes into the training set while the minor chunk forms the test set. The tree construction involves the following steps:

1. If the training set has  $N$  cases, a sample of  $N$  cases is randomly selected from the original dataset. This sample corresponds to training set that is used for growing the tree.
2.  $m$  variables out of the  $M$  input variables are chosen randomly, the node is split based on the best split on this  $m$  value.  $m$  is held constant while growing the forest.
3. Each tree in the forest is grown to the largest possible extent. No Trimming or Pruning is performed on the individual trees.
4. All classification trees thus formed are combined to form the random forest. Since it can fix problem of overfitting on large dataset and can train/test rapidly on complex data set, it is sometimes referred to as *Operational Data mining* technique.

Each classification tree is exclusive and is voted for a class. Finally, a solution class is constructed based on the maximum votes assigned.

#### d) Association Rule Mining (ARM)

Association-rule mining discovers interesting relations between a set of attributes in datasets[16]. The datasets and their inter-relationship can be represented as association rules. This information can be used for making strategic decisions about different activities such as, promotional pricing, shelf management and so on[17]. Traditional Association rule mining involves a data analyst being given datasets of different companies for the purpose of discovering patterns or association rules that exist between the datasets[18]. Although, we can achieve sophisticated analysis on these extremely large datasets in a cost-effective manner[19], it poses security risk[20] for the data owner whose sensitive information can be deduced by the dataminer[21]. Even today, association rule mining is one of the widely used pattern discovery methods in KDD.

Solving an ARM problem basically involves traversing the items in a database, which can be done using various algorithms based on the requirement[22]. ARM algorithms are primarily categorised into BFS (Breadth First Search) and DFS (Depth First Search) methods based on the strategy used to traverse the search space[23]. The BFS and DFS methods are further classified into Counting and Intersecting, based on how the support values for the itemsets are determined. The algorithms Apriori, Apriori-TID and Apriori-DIC are based on BFS with Counting strategies, while the Partition algorithm is based on BFS with Intersecting strategies. The FP-Growth algorithm on the otherhand, is based on DFS with Counting strategies while ECLAT is based on DFS with Intersecting[24][25]. These algorithms can be optimized specifically for improving the speedup [26][27].

*BFS with Counting Occurences:* The common algorithm in this category is the Apriori algorithm. It utilizes the

downward closure property of an itemset, by pruning the candidates with infrequent subsets before counting their supports. The two metrics to be considered while evaluating the association rules are: *support* and *confidence*. BFS offers the desired optimization by knowing the support values of all subsets of the candidates in advance. The limitation of this approach is increased computational complexity in rule extraction from a large database. Fast Distributed Mining(FDM) algorithm is a modified, distributed and unsecured version of the Apriori algorithm[28]. The advancements in data mining techniques, have enabled organizations in using data more efficiently.

In Apriori, the candidates of a cardinality  $k$  are counted by a single scan of the entire database. Looking up for the candidates in each transaction forms the most crucial part of the Apriori Algorithm. For this purpose, a hashtree structure is used[29]. Apriori-TID an extension of Apriori, represents each transaction based on the current candidates it contains, unlike normal Apriori that relies on raw database. Apriori-Hybrid combines the benefits of both Apriori and Apriori-TID. Apriori-DIC another variation of Apriori, tries to soften the separation that exists between the processes, counting and candidate generation. This is done by using a prefix-tree.

*BFS with Intersections:* A Partition Algorithm is similar to the Apriori algorithm that uses intersections rather than counting occurrences for the determination of support values. The partitioning of itemsets could result in the exponential growth of intermediate results beyond the physical memory limitations. This problem can be overcome, by splitting the database up into a number of chunks that are smaller in size and each chunk is treated independently. The size of a chunk is determined such that all intermediate lists can fit into memory. An additional scan can optionally be performed to ensure that the itemsets are not only locally frequent but also are globally frequent.

*DFS with Counting Occurences:* In Counting, a database scan for each reasonable sized candidate set is performed. Because of the involvement of computational overhead in database scanning, the simple combination of DFS and Counting Occurences is practically irrelevant. FP-Growth on the otherhand uses a highly compressed representation of transaction data called *FP-Tree*. An FP-Tree is generated by counting occurrences and performing DFS.

*DFS with Intersections:* The algorithm *ECLAT* combines DFS with the list intersections to select agreeable values. It makes use of an optimization technique called *Fast Intersections*. It does not involve the process of splitting up of the database since complete path of classes beginning from the root would be maintained in the memory. As this method eliminates most of the computational overhead the process of mining association rules becomes faster.

### III. CLUSTERING

Clustering is one of the widely used discovery methods in data mining. It allows to group a set of data in such a way that, Intra-Cluster similarity are maximized while minimizing the Inter-Cluster similarity are minimized. Clustering involves unsupervised learning of a number of classes that are not known in advance. The clustering algorithms can be broadly clasified into the following types and are listed in Fig.3

1. Connection Based or Hierarchical Clustering
2. Centroid Based
3. Distribution Based
4. Density Based
5. Recent Clustering Techniques and
6. Other Clustering Techniques

a) *Connection Based Clustering*

Connection Based (Hierarchical) clustering, is based on the idea of objects being more related to closer objects than to the distant objects. The Connection Based Clustering algorithms consider the distance between the objects to connect them to form *clusters*. These algorithms provide an extensive hierarchy of merging clusters at particular distances, instead of single partitioning of dataset. A *Dendrogram* is used to represent clusters. Its *y-axis* shows the merging distance of the clusters

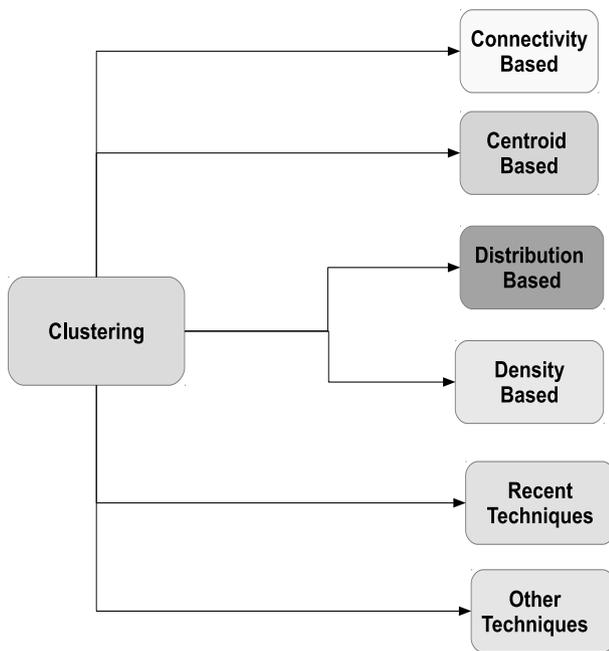


Figure 3: Types of Clustering

and the *x-axis*, for placing the objects, ensuring that the clusters do not mix. There are various types of Connection based clustering based on the way the distances are computed such as: Single-Linkage Clustering that involves determining of the minimum of object distances, Complete-Linkage Clustering where the maximum of object distances is computed and

Unweighted Pair Group Method with Arithmetic Mean (UPGMA), or Average Linkage Clustering. Selecting appropriate clusters from the available hierarchy of clusters, could be achieved either using Agglomerative or Divisive Clustering. In Agglomerative Clustering, we begin with single objects and conglomerate them into clusters while in Divisive clustering, we start with the complete data set and isolate it into segments.

b) *Centroid Based Clustering*

Centroid-based clustering may have clusters that are represented by a vector, which necessarily is not a member of the data set or may have clusters strictly restricted to the members of the dataset. In *k-means Clustering* algorithm, the number of clusters is limited to size *k*, it is required to determine *k* cluster centers and assigning objects to their nearest centers.

The algorithm is run multiple times with different *k* random initializations to choose the best of multiple runs[30]. In *kmedoid clustering*, the clusters are strictly restricted to the members of the dataset while in *k-medians clustering*, only the medians are chosen to form a cluster. The main disadvantage of these techniques is that the number of clusters *k* is selected beforehand. Furthermore, they result in incorrectly cut borders in between the clusters.

c) *Distribution Based Clustering*

Distribution-based clustering technique forms clusters by choosing objects that belong more likely to the same distribution. One of the most commonly preferred distribution techniques is the Gaussian Distribution. It suffers from the overfitting problem where a model cannot fit into set of training data.

d) *Density Based Clustering*

In this type of clustering, an area that is having higher density than the rest of the data set is considered as a cluster. Objects in the sparse areas are considered to be noise and border points. There are three commonly used Density-based Clustering techniques namely: DBSCAN, OPTICS and Mean- Shift. DBSCAN is based on connecting points that satisfy a density criterion within certain distance thresholds. The cluster thus formed may consist of all density-connected objects and objects that are within these objects range free to have an arbitrary shape.

e) *Recent Clustering Techniques*

All the standard clustering techniques fail for highdimensional data and so some of the new techniques are being explored. These techniques fall into two categories namely: Subspace Clustering and Correlation Clustering. In Subspace Clustering, the clustering model specifies a small list of attributes that should be considered for the formation of a cluster while in Correlaton Clustering, the model along with this list of attributes it also provides the correlation between the chosen attributes.

f) Other Techniques

One of the most basic clustering techniques is the BSAS(Basic Sequential Algorithmic Scheme). Given the distance  $d(p, C)$  between a vector point  $p$  and a cluster  $C$ , the maximum number of clusters allowed  $q$  and threshold of dissimilarity  $\theta$ , the BSAS constructs the clusters even when the number of clusters to be formed is not known in advance.

Every newly presented vector is either assigned to an already existing cluster or a new cluster is created, depending on the distance to the already present clusters.

g) Clustering applications in IDS

Clustering technique may be effectively used in the process of Intrusion Detection. The setup is depicted

in Fig. 4. Alerts generated by multiple IDSs belonging to both Network and Host types are logged into a centralized database. The alert messages arriving from different IDSs will be in different formats. Before passing them into the server, a preprocessing step is needed to bring them all into some uniform format [31].

Best effort values are chosen for the missing attributes during the preprocessing stage. The timestamp information may have to be converted into seconds for the sake of comparison. Different IDSs may use different conventions for naming a single event and hence it is required to standardize

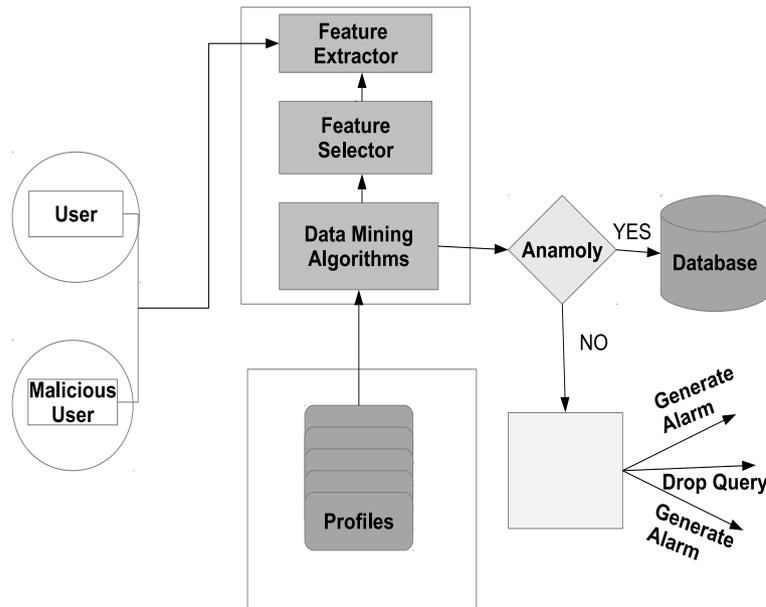


Figure 4: Use of Clustering in IDS

the messages. Each alert may be added with a unique ID to keep track of the alerts. After preprocessing and normalizing alerts they are passed to the first phase to perform filtering and labeling functions. To minimize the number of Alerts, it is a good idea to employ Alert Fusion during which alerts with same attributes that differ by a small amount of time are fused together. Alert Fusion makes the generalization process fast. Generalization involves the addition of hierarchical background knowledge into each attribute. On every iteration of this process, the selected attribute is generalized to the next higher level of hierarchy and those alerts which have become similar by now are grouped together.

perform all data mining operations efficiently. The two types of privacy concerned data mining techniques are:

1. Data privacy
2. Information privacy

Data privacy focuses on the modification of the database for the protection of sensitive data of the individuals while Information privacy focuses on the modification for the protection of sensitive knowledge that can be deduced from the database.

Alternatively we can say that Data privacy is concerned about providing privacy to the input while Information privacy on the otherhand is about providing privacy to the output. Preserving personal information from revelation is the main focus of a PPDM algorithm[32]. The PPDM algorithms rely on analysing the mining algorithms for any side effects that are acquired during Data privacy. The objective of Privacy Preserving Data Mining is building algorithms that transform the original data in some mannner, so that both the private data and knowledge are not revealed even after a successful mining process. Only when some relevant adequate benefit is found resulting from the access, the privacy laws would allow the access.

#### IV. PRIVACY PRESERVING DATA MINING (PPDM)

Privacy Preserving Data Mining techniques aim at the extraction of relevant knowledge from large volumes of data while protecting any sensitive information present in it. It ensures the protection of sensitive data to conserve privacy and still allowing us to

Multiple parties may sometimes wish to share private data resulting after a successful aggregation[33] without disclosing any sensitive information from their end[34]. Consider for example, different Book stores with respective sales data that is in a way considered to be highly sensitive, may wish to exchange partial information among themselves to arrive at the aggregate trends without disclosing their individual store trends. This requires the use of secure protocols for sharing the information across multiple parties. Privacy in such cases should be achieved with high levels of accuracy[35].

The data mining technology by principle is neutral in terms of privacy[36]. The motive for which a data mining algorithm is used could either be good or malicious[37]. Data mining has expanded the investigation possibilities[38] to enable researchers to exploit immense datasets on one hand[39], while the malicious use of these techniques on the other hand has introduced threats of serious nature against protection of privacy[40].

Discovering the base of privacy preserving data mining

Table 1: Research Progress in PPDM

Authors	Algorithm	Performance	Future enhancement
Boutet et al.(2015)[45]	kNN	Better than Randomization scheme	Can consider all attacking models
Tianqing et al.(2015)[46]	Correlated Differential Privacy (CDP)	Enhances the utility while answering a large group of queries on correlated datasets	Can be experimented with Complex Applications
Bharath et al.(2015)[47]	PP k-NN classifier	Irrespective of the values of k, it is observed that SRkNN is around 33% faster than SRkNN. E.g., when k=10, the computation costs of SRkNN and SRkNN are 84.47 and 127.72 minutes, respectively (boosting the online running time of Stage 1 by 33.86%)	Parallelization is not used
Nethravathi et al.(2015)[48]	PPDM	Reduced misplacement clustering error and removal of data that is sensitive and correlated	Works only for numerical data
Mohammed et al.(2014)[49]	Differential Privacy	More secured under the Semi-Honest model	Overcoming Privacy Attack
Vaidya et al.(2014)[50]	Distributed RDT	Lower Computation and Communication cost	Limited information that is still revealed must be checked
Lee(2014)[51]	Perturbation methods	Capable of performing RFM Analysis	Partial disclosure is still possible

algorithms and connected privacy techniques is the need of the hour[41]. We are required to answer few questions in this regard such as

1. Evaluation of these algorithms with respect to one another
2. Should privacy preserving techniques be applied to each of the data mining algorithms? Or for all applications?
3. Expanding the places of usage of these techniques.
4. Investigating their use in the fields of Defense and Intelligence, Inspection and Geo-Spatial applications.
5. The techniques of combining confidentiality, privacy and trust with high opinion to data mining.

To answer these questions, research progresses in both data mining and privacy are required. Proper planning towards developing flexible systems is essential[42]. Few applications may demand *pure data mining* techniques while few others may demand *privacy-preserving data mining*[43]. Hence we require flexible techniques in data mining that can cater to the the changing needs[44]. The research progress made so far in the area of PPDM is listed in Table 1.

*Distributed Privacy Preserving Data Mining(DPPDM):*

The tremendous growth of internet in the recent times is creating new opportunities for distributed data mining[52], in which, mining operations performed jointly using their private inputs[53]. Often occurrence of mining operations between untrusted parties or competitors, result in privacy leakage[54]. Thus, Distributed Privacy Preserving Data Mining(DPPDM)[10][55] algorithms require a high level of collaboration between parties to deduce the results or to share mining results that are not sensitive. This could sometimes result in the disclosure of sensitive information.

Distributed data mining are classified as Horizontally Partitioned Data and Vertically Partitioned Data. In a Horizontally partitioned data framework, each site maintains complete information on an unique set of entities, and the integrated dataset consists of the union of all of these datasets. Vertically Partitioned Data framework on the otherhand involves each site, maintaining different types of information and each dataset and has only limited information about same set of entities.

Privacy feature can limit the information leakage caused by the distributed computation techniques[56].

Each non-trusting party can compute its own functions for unique set of inputs, revealing only the defined outputs of the functions. Apart from hiding sensitive information, the privacy service also controls the information and its uses by involving various number of negotiations and tradeoffs between hiding and sharing.

All efficient PPDM algorithms are based on the assumption that it is acceptable to release the intermediate results obtained during the data mining operations. Encryption techniques solve the data privacy problem and their use would make it easy to perform data mining tasks among mutual untrustworthy parties, or between competitors. Due to its privacy concern, Distributed Data Mining Algorithms employ encryption techniques. Encryption is used in both approaches (horizontally and vertically partitioned data) of Distributed Data mining without much stress on the efficiency of encryption technique used.

If the data are stored on different machines and partitioning is done row-wise, it is called horizontal partitioning and if the data are stored and partitioned column wise then it is called vertical partitioning. An overview of the same is depicted in Fig.5.

The objective of data mining techniques is to generate high level rules or summaries and generalize across populations, rather than revealing information about individuals but they work by evaluating individual data that is subject to privacy concerns. Since much of this information held by various organizations has already been collected, providing privacy is a big challenge. To prevent any correlation of this information,

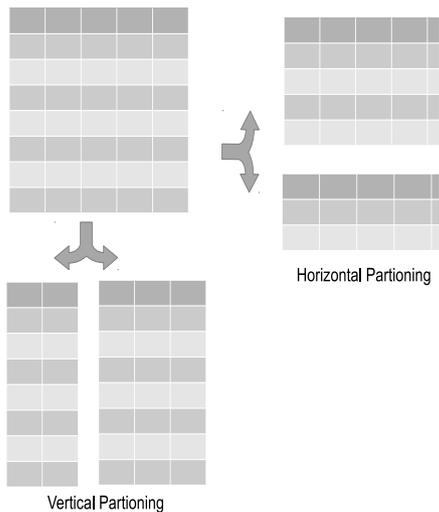


Figure 5: Horizontal and Vertical Partitioning Techniques

control and individual safeguards must be separated to be able to provide acceptable privacy. Unfortunately, this separation makes it difficult to use the information for the identification of criminal activities and other purposes that would benefit the society. Proposals to share information across agencies to combat terrorism and other criminal activities, would also remove the safeguards imposed by separation.

Many of the complex socio-technical systems suffer from an inadequate risk model that focuses on the use of Fair Information Practice Principles (FIPPs). Anonymization suffers from the risk of failure, since the circumstances surrounding its selection are ignored. A Hybrid approach that combines privacy risk model with an integrated anonymization framework involving anonymization as the primary privacy risk control measure can be considered instead [57].

**Public-Key Program Obfuscation:** The process of making a program uncomprehensible without altering its functionality is called Program Obfuscation. A program that is obfuscated should be a *virtual black box* meaning, if it is possible for one to compute something from it, it should also be possible to compute the same even from the input-output behavior of the program. Single-Database Private Information Retrieval can be considered a type of public-key program obfuscation. Given a program  $p$  from a class of programs  $C$ , and a security parameter  $s$ , a public-key program obfuscation function compiles  $p$  into  $(P, Dec)$ , where  $P$  on any input computes an encryption of what  $p$  would compute on the same input and the decryption algorithm  $Dec$  decrypts the output of  $P$ . That is, for any input  $i$ ,  $Dec(P(i)) = p(i)$ , but for given code  $P$  it is impossible to distinguish which  $p$  from the class  $C$  was used to produce  $P$ . The program encoding length  $|P|$  must depend only on  $|p|$  and  $s$ , and the output length of  $|P(i)|$  must polynomially depend only on  $|p(i)|$  and  $k$ .

**Secure Multi-party Computation:** Distributed computing involves a number of distinct, and connected computing devices that wish to carry out a combined computation of some function. For example, servers holding a distributed database system, may wish to update their database. The objective of secure multiparty computation is to allow parties to carry out distributed computing tasks in a secure way [33]. It typically involves the parties carrying out a computation based on their private inputs and neither of them willing to disclose its own input to other parties. The problem is conducting such a computation by preserving the privacy of their inputs. This problem is called the Secure Multi-party Computation problem (SMC) [34]. Consider the problem of two-parties who wish to securely compute the median. The two parties have with them two separate input sets  $X$  and  $Y$ . The parties are required to jointly compute the median of the union of their sets  $X \cup Y$ , without revealing anything about each other's set. Association Rules can be computed in an environment where different information holders have different types of information about a common set of entities.

## V. INTRUSION DETECTION SYSTEM (IDS)

Intrusion detection systems aim at the detection of an intrusion on its onset [58]. A high level of human expertise and significant amount of time are required for

the development of a comprehensive IDS[59]. However, IDSs that are based on the Data Mining techniques require less expertise and yet they perform better. An Intrusion Detection System detects network attacks against services that are vulnerable [60], attacks that are data driven on applications, privilege escalation[61], logins that are un-authorized and access to files that are sensitive in nature[62]. The data mining process also efficiently detects malware from the code[63], which can be used as a tool for cyber security[64][65]. An overview of an Intrusion Detection System is presented in Fig 6.

An IDS is basically composed of several components such as, sensors, a console monitor and a central engine[66]. Sensors generate security events while all events and alerts are monitored and controlled by the Console Monitor and the Central Engine records events in a database and generate alerts based on a set

of rules[67]. An Intrusion detection system[68] can be classified depending on the location and the type of Sensors and based on the technique used by the Central engine for the generation of alerts. A majority of IDS implementations, involve all of the three components integrated into a single device.

Current virus scanner methodology makes use of two parts namely a Detector based on signatures and a Classifier based on the heuristic rules for the detection of new viruses. The signature-based detection algorithms rely on signatures that are unique strings of known malicious executables for the generation of detection models. The disadvantages of this approach are: it is more time-consuming and fails in detecting new malicious executables. Heuristic classifiers on the other hand are generated by a set of virus experts for the detection of new malicious executables.

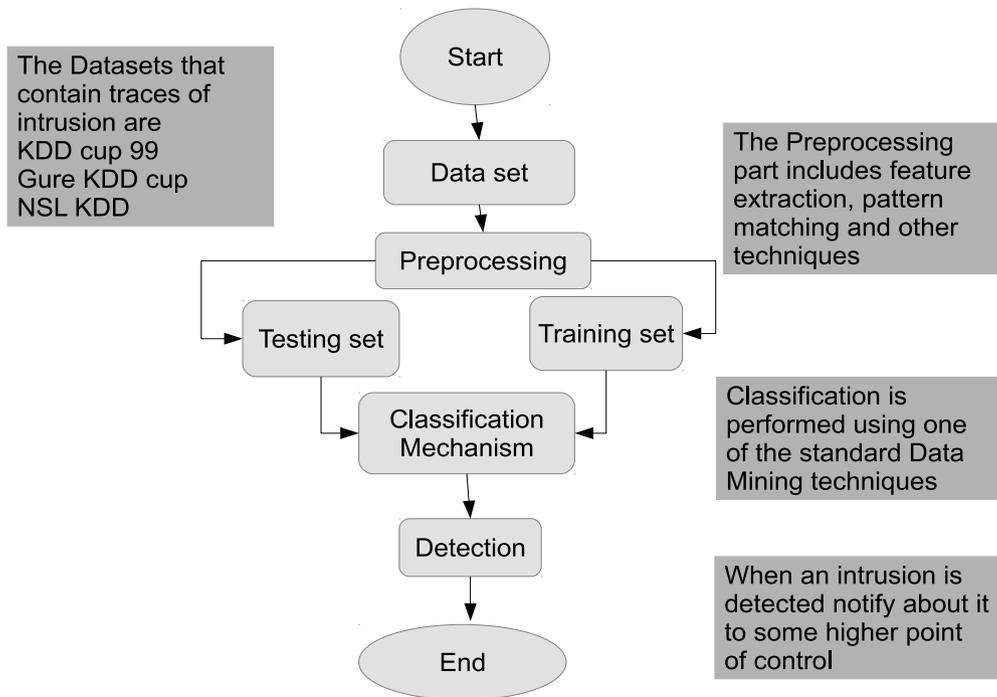


Figure 6: An overview of an Intrusion Detection System

a) Types of IDS

An intrusion could be detected either on a individual system or on a network and accordingly we have three types of IDS namely: Network Based, Host Based and Hybrid IDS.

i. Network Based IDS

Because of their increasingly vital roles in modern societies, computer networks have been targeted by enemies and criminals. For the protection of our systems, it is very essential to find the best possible solutions. Intrusion prevention techniques such as, authentication techniques involving passwords or biometrics[69], programming errors avoidance, and

protection of information using encryption techniques have been widely used as a first line of defense. Intrusion prevention techniques as the sole defense mechanism are not sufficient enough to combat attacks. Hence, it can therefore be used only as a second line of defense for the protection of computer systems[70].

An Intrusion Detection system must protect resources such as accounts of users[71], their file systems and the system kernels of a target system and must be able enough to characterize the legitimate or normal behavior of these resources by involving techniques that compare the ongoing system activities with already established models and to identify those activities that are intrusive[72][73]. Network packets are

the data source for Network-Based Intrusion Detection Systems. The *NIDS* makes use of a network adapter to listen to and analyse network traffic as the packets travel across the network. A Network based IDS generates alerts upon detecting an intrusion from outside the perimeter of its enterprise[74]. The network based IDSs are categorically placed at strategic points on LAN to observe both inbound and outbound packet[75]. Network based IDSs are placed next to the firewalls to alert about the inbound packets that may bypass the firewall[76]. Few Network-Based IDSs take custom signatures from the user security policy as input, permitting limited detection of security policy violations[77]. When packets that contain intrusion originated from authorized users, the IDS may not be able to detect[78][79].

#### *Advantages*

Some of the advantages of a Network Based IDS are as follows:

1. For enhanced security against attacks, they can be made invisible.
2. Are capable of monitoring larger networks.
3. They can function without interfering with the normal operation of a network[80].
4. It is easy to fit in an IDS into an existing network.

#### *Disadvantages*

The disadvantages are as follows:

1. Not capable enough to analyze encrypted information coming from virtual private networks.
2. Their success most of the times depend on the capabilities of the intermediate switches present in the network.
3. When the attackers fragment their packets and release them, the IDS could become unstable and crash.

#### ii. *Host Based IDS*

In a Host-based IDS, the monitoring sensors are placed on network resources nodes so as to monitor logs that are generated by the Host Operating System or application programs.

These Audit logs contain records of events or activities that are occurring at individual Network resources[81]. Since a Host- Based IDS is capable of detecting attacks that cannot be seen by a Network based IDS, an attacker can misuse one of trusted insiders[82]. A Host based system utilizes Signature Rule Base that is derived from security policy that is specific to a site. A Host Based IDS can overcome all the problems associated with a Network based IDS as it can alarm the security personnel with the location details of intrusion, he can take immediate action to thwart the intrusion. A Host based IDS can also monitor any unsuccessful attempts of an attacker. It can also maintain separate records of user login and user logoff actions for the generation of audit records.

#### *Advantages*

Some of the advantages of a Host Based IDS are as follows:

1. Can detect attacks that are not detected by a Network Based IDS.
2. Operates on Operating System audit log trails, for the detection of attacks involving software integrity breaches.

#### *Disadvantages*

The disadvantages are:

1. Certain types of DoS(Denial of Service)attacks can disable them[83].
2. Not suited for detecting attacks that target the network.
3. Difficult to configure and manage every individual system.

#### iii. *Hybrid IDS*

Since Network and Host-based IDSs have strengths and benefits that are unique over one another, it is a good idea to combine both of these strategies into the next generation IDSs[84]. Such a combination is often referred to as a Hybrid IDS. Addition of these two components would greatly enhance resistance to few more attacks.

#### a. *DM techniques for IDS*

Some of the techniques and applications of data mining required for IDS include the following

1. Pattern Matching
2. Classification and
3. Feature Selection

*Pattern Matching:* Pattern Matching is a process of finding a particular sequence of a part of data (substring or a binary pattern), in the whole data or a packet to get a desired information[87]. Though it is fairly rigid, it is indeed simple to use. A Network Based IDS succeeds in detecting an intrusion only when the packet in question is associated with a particular service or, destined to or from a particular port. That is, only few fields of the packet such as Service, Source/Destination port address and few others have to be examined thereby reducing the amount of inspection to be done on each packet.

However, it makes it difficult for systems to deal with Trojans and their associated traffic that can be moved at will. The pattern matching can be classified into two categories based on the frequency of occurrence namely:

- a) Frequent Pattern Matching and
  - b) Outlier Pattern Matching
- a) *Frequent Pattern Matching*

These are the type of patterns which occur frequently in an audit data, i.e., the frequency of occurrence of these patterns is more compared to other patterns in the same data[82].

Determining frequent patterns in a big data helps in analyzing and forecasting of a particular characteristic of the data. For example, by analyzing the sales information of an organization, frequent pattern matching might help to predict the possible sales outcome for the future. It also helps in decision making. The frequent pattern mining in ADAM project data is done by mining the repository for attack-free (train) data which is compared with the patterns of normal profile (train) data. A classifier is used to reduce the false positives.

b) *Outlier Pattern Matching*

Patterns that are unusual and are different from the remaining patterns and that are not noise are referred to as Outlier Patterns. Preprocessing phase eliminates noise as it is not a part of the actual data while outliers on the other hand cannot be eliminated. Outliers exhibit deviating characteristics as compared to the majority of other instances. Outliers patterns are not

usual and they occur less frequently and for this reason will have minimal support in the data. These patterns can quite often point out some sort of discrepancy in data such as transactions that are fraudulent, intrusion, abnormal behavior, economy recession etc.,. The outlier pattern mining algorithms can be of two types, one that looks for patterns only at fixed time intervals, and the other that calculates monitors patterns at all times. Outlier pappers make use of special data structures such as Suffix Tree and other String Matching Algorithms.

*Classification:* Classification makes use of training examples for learning a model and to classify samples of data into known classes[88]. A wide range of classification techniques ranging from Neural Networks, Decision Trees, Bayesian classifier[89], Bayesian Belief Networks and others are used in applications that involve Data Mining techniques. Classification typically involves steps that are outlined below:

Table 2: Research Progress in IDS

Authors	Algorithm	Performance	Future enhancement
M Vittapu et al.(2015)[85]	SVM Classification	TPR of 96% and FPR of 5%	Can be experimented with other techniques
Mitchell et al.(2015)[61]	Behavior Rule Analysis	Better performance	Can be tested with other techniques
Jabez J et al.(2014)[98]	Hyperboli Hopfiel Neural Network(HHNN)	Detection rate of about 90%	Can be improved
S Abadeh et al.(2014)[151]	Genetic Fuzzy System	Best tradeoff in terms of the mean F-measure,the average accuracy and the false alarm rate	A Multi-objective Evolutionary Algorithm for maximizing performance metrics may be considered
Soni et al.(2014)[86]	Feature Selection	Better classification	Can consider NSL-KDD

1. Creation of a training dataset
2. Identification of classes and attributes
3. Identification of attributes that are useful for classification
4. Relevance analysis
5. Learning the Model using training examples
6. Training the set
7. Using the model for the classification of unknown data samples.

*Bayesian Classifiers:* The Naive Bayesian approach assumes, the attributes to be independent in condition. Although it works under this assumption, the Naive Bayesian classifiers yield results that are satisfactory because they focus on identifying the classes for the instances instead of their probabilities. Spam Mail classification and Text classification applications extensively use Naive Bayesian classifiers for they are less error prone. However, their disadvantage is that they require probabilities in advance. The probability information that is required by them is extremely huge which consist number of classes, their attributes and the maximum cardinality of attributes. The space and computational complexity of these classifiers increase exponentially.

*Support Vector Machine(SVM):* Support Vector Machine is one of the learning methods extensively used for the Classification and Regression analysis of Linear and Non-linear data[90]. It maps input feature vectors into a higher dimensional space using non-linear mapping techniques. In SVM, the classifier is created by the linear separation of hyperpalnes and linear separation is achieved using a function called kernel.The Kernel transforms a linear problem by mapping it into feature spaces.

Some of the commonly used kernel functions are Radial basis, sigmoid neural nets and polynomials. Users specify one of these functions while training the classifier and it selects support vectors along the surface of this function. The SVM implementation tries to achieve maximum separation between the classes[91]. Intrusion detection system involves two phases namely training and testing. SVMs are capable of learning a larger set of patterns and can provide better classification, because the categorizing complexity is independent of the feature space dimensionality[92]. SVMs can update the training patterns dynamically with the availability of new pattern during classification. For the efficient classification it is required to reduce the

dimensionality of the dataset. To do this we have *Feature Selection*.

### iii. *Feature Selection(FS)*

The process of reducing the dataset dimensionality by selecting a subset of the features from the given set of features is called Feature Selection[93]. FS involves discarding of redundant and irrelevant features. FS is considered to be an efficient machine learning technique that helps in building classification systems which are efficient. With the reduction in subset dimensionality, the time complexity is reduced with improved accuracy, of a classifier. Information Gain is a proposition of feature selection that can be used to compute entropy cost of each attribute. An entropy cost can be called as a rank. Rank of each feature represents its importance or association with an solution class that is used to recognize the data. So a feature with comparatively higher rank will be one of the most important features for classification. The three standard approaches that are commonly followed for feature selection are embedded technique, filter technique, and wrapper technique.

FS runs as a part of data mining algorithms, in Embedded technique. Feature selection is independent of the classifier used in case of Filter method, while in Wrapper method features are chosen specifically to the intended classifier. Filter method uses an arbitrary statistical way for the selection of features whereas wrapper method uses a learning algorithm to find the best subset of features. Wrapper approach is more expensive and requires more computational time than the filter approach but gives more accurate results compared to filter technique.

## VI. PHISHING WEBSITES CLASSIFICATION

In the art of emulating a website of a trusted and creditable firm with the intention of grabbing users' private information (username, password) is called *phishing*. Fake websites are usually created by dishonest people to masquerade honest websites. Users unknowingly lose money due to phishing activities of attackers. Online trading therefore demands protection from these attacks and is considered a critical step. The prediction and classification accuracy of a website depends on the goodness of the extracted features. Most of the internetusers feel safe against phishing attacks by utilizing antiphishing tool, and hence the anti-phishing tools are required to be accurate in predicting phishing[94]. Phishing websites give us a set of clues within its content parts and through security indicators of the browsers[95]. A variety of solutions have been proposed to tackle the problem of phishing. Data mining techniques involving Rule based classification[96] serve as promising methods in the prediction of phishing attacks.

Phishing attack typically starts by, attacker sending an email to victims requesting personal information to be disclosed, by visiting a particular

URL[97]. Phishers use a set of mutual features to create phishing websites to carry out proper deception[98]. We can exploit this information to successfully distinguish between phishy and non-phishy websites based on the extracted features of the website visited[94]. The two approaches that are commonly used in the identification of phishing sites are: black-list based, which involves comparison of the requested URL with those that are present in that list and Heuristic based method that involves the collection of certain features from the website to label it either as phishy or legitimate[99]. The disadvantage of Black-list based approach is that the black-list can not contain all phishing websites since, a new malicious website is launched every second[100]. In contrast, a Heuristic-based approach can recognize fraudulent websites that are new[101]. The success of Heuristic-based methods depend on the selection of features and the way they are processed. Data mining can be effectively used here to find patterns as well as relations among them[102]. Data mining is considered to be important for taking decisions, since decisions are made based on the patterns and rules derived using the data mining algorithms[103].

Although there is substantial progress made in the development of prevention techniques, phishing still remains a threat since most of the counter measures techniques in use are based still on reactive URL black-listing[104]. Since Phishing Web sites will have shorter life time these methods are considered to be inefficient. Newer approaches such as Associative Classification (AC) are more suitable for these kinds of applications. Associative Classification technique is a new technique derived by combining Association rule and Classification techniques of data mining[105]. AC typically includes two phases; the training phase to induce hidden knowledge (rules) using Association rule and the Classification phase to construct a Classifier after pruning useless and redundant rules. Many research studies have revealed that AC usually shows better classifiers with reference to error rate than other standard classification approaches such as decision tree and rule induction.

## VII. ARTIFICIAL NEURAL NETWORKS(ANN)

An Artificial Neural Network is basically a connected set of processing units. Each connection has a specific weight that determines how one unit affects the other. Few of these units act as input nodes and few other as output nodes and remaining nodes consists of hidden layer. Neural network performs functionally, a mapping from input values to output values by activating each input node and allowing it to spread through the hidden layer nodes to the output nodes. The mapping is stored in terms of weight over connection. Fig. 7 shows the structure of HHNN[62].

ANN is one of the widely used techniques in the field of intrusion detection. ANN techniques are classified into three categories namely:

1. Supervised Intrusion Detection,
2. Unsupervised Intrusion Detection, and
3. Hybrid Intrusion detection.

Supervised Intrusion Detection based on ANN includes Multi Layer Feed Forward (MLFF) Neural Networks and Recurrent Neural Networks. Since, the number of training sets is huge and their distribution is imbalanced, the MLFF neural networks can easily reach the local minimum and hence have lower stability. The precision rate of a MLFF neural network is low for less frequent attacks. Supervised IDS exhibits lower detection performance than SVM and Multivariate Adaptive Regression Splines(MARS). Unsupervised Intrusion Detection based on ANN classifies test data and separates normal behaviors from abnormal ones. Since it does not need retraining, it can greatly improve the analysis of new data. The performance of Unsupervised ANN is also lower for low frequent attacks achieving a lower detection precision. Hybrid ways of combining supervised ANN and unsupervised ANN and combining ANN with other data mining techniques for the detection of intrusion can be achieved to overcome the limitations of the basic types of ANN. A hybrid approach involving SOM and Radial Basis Function(RBF) networks is comparatively more efficient than Intrusion Detection based on RBF networks alone. A hybrid model that uses a combination of Flexible Neural Tree, Evolutionary Algorithm and Particle Swarm Optimization (PSO) is highly efficient. Hybrid ANN that uses a combination of Fuzzy Clustering technique with

ANN reduces the training set into subsets that are smaller in size, thereby improving the stability of individual ANN for low-frequent attacks. So we can say that for intrusion detection based on ANNs, hybrid ANN has been the trend. Different ways of constructing Hybrid ANN influences the performance of intrusion detection. Hence it is required to construct different Hybrid ANN models to serve different goals. There are various Hybrid approaches being utilized for intrusion detections and one such model is the Hyperbolic Hopfield Neural Network(HHNN). Anomaly detection assumes that the intrusions always return as a number of deviations from the normal patterns. HHNN technique studies the relationship between the two sets of information, and generalizes it in getting new input-output pairs reasonably. Neural networks can be used hypothetically for the identification of attacks and look for these attacks in the audit stream. Since there is no reliable method at present to realize causes of association, it cannot clarify the reason behind the classification of the attack. The research progress made in HHNN is summarized in Table 3.

### VIII. ANOMALY DETECTION/OUTLIER DETECTION

Anomaly detection is a process that involves finding nonconforming patterns to the expected behavior. Such patterns are called *anomalies*. Different application domains term them differently as outliers or aberration or surprises or peculiarities or

Table 3: Research Progress in ANN

Authors	Algorithm	Performance	Future enhancement
C Cortes et al.(2016)[106]	Theoretical framework for analyzing and learning artificial neural networks	Optimizes generalization performance	Can be applied for different optimization techniques and network architectures.
D T Bui et al.(2015)[107]	ROC and Kappa Index	MLP (90.2 %), SVM (88.7 %), KLR (87.9 %), RBF (87.1 %) and LMT (86.1 %).	Information Gain Ratio as feature selection can be tried.

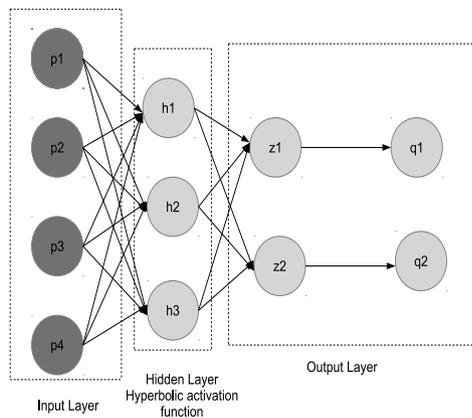


Figure 7: An overview of a HHNN

contaminants. Anomalies and Outliers are the two commonly used terms in this context. Anomaly detection applications are fraud detection of credit or debit cards, health care and insurance. It is also used for intrusion detection and fault detection in a safety critical system and for the detection of enemy activities.

Anomalous patterns mostly deviate from the normal patterns. The figure shown in Fig. 8 plots anomalies on a two dimensional data set. The regions *N1* and *N2* are considered normal, because majority of the observations lie in these regions. Points *O1* and *O2* that are far away from these regions and points in region *O3*, are considered anomalies. Although anomalies are induced in the data for a number of reasons, all of them have the common characteristic that they are interesting to analyze. This interestingness of outliers is a prime feature of anomaly or outlier detection. Anomaly detection is related to, but not as same as noise removal and noise accommodation, but it must ceratinly deal with unwanted noise in the data. Noise is an unwanted part to the analyst and acts as a hurdle to data analysis. Noise removal is therefore necessary, since unwanted data must be removed before performing the data analysis. Novelty Detection is a topic related to anomaly detection, which detects any previously unidentified novel patterns in the data. The detected novel patterns are incorporated into the normal model, that makes it different from Anomaly Detection. Different solutions that exist for anomaly detection will also work for novel Detection and *vice versa*. Hence in Anomaly Detection a region is defined, where the observations conforming to the region are considered normal and the non-conforming observations are considered anomolous.

a) Challenges

Some of the challenges the researchers face with respect to Anomaly Detection are:

1. Defining a normal region, where all normal behaviors exist is difficult. The boundary between normal and anomalous behavior has a very thin differentiation, meaning that an observation that lies closer to the boundary could be normal, and *vice versa*.
2. When the attackers masquerade to make the anomalous observations to appear normal, defining normal behavior becomes complicated.
3. Normal behaviors evolve and what is currently considered as normal might not be the same in the future.
4. Different application domains have different notions of anomaly. For instance fluctuations in body temperature marks an anomaly in the medical domain, while the fluctuations in marketing domain might be considered as normal. Therefore the application of a technique developed, cannot be generic.

5. Labeled data used by anomaly detection techniques for training/validation of models is not available freely. It is challenging to distinguish and remove noise from the data. The anomaly detection issue, is therefore hard to tackle with. Most of the anomaly detection techniques that exist, can only solve a problem formulation that is domain specific and is induced by factors such as category of the data, labeled data availability, anomaly types, and so on.

b) Data Mining Mechanisms for Anomaly Detection

An Intrusion Detection System can generally be implemented using the following two techniques:

1. Signature Based IDS and
2. Anomaly Based IDS.

Signature Based IDS makes use of Attack signatures that are explicitly defined and detect intrusions by Blacklisting.

It is ineffective against new types of attacks which makes it susceptible to evasion methods.

Anomaly Based IDS on the other hand, records normal behavior and classifies the deviations from normal behavior as anomalies. It is considered to be robust and reliable to unknown attacks and prevent attacks from malicious users who improvise their attacking strategy. The widely used implementation of Anomaly Based IDS is by the extensive use of data

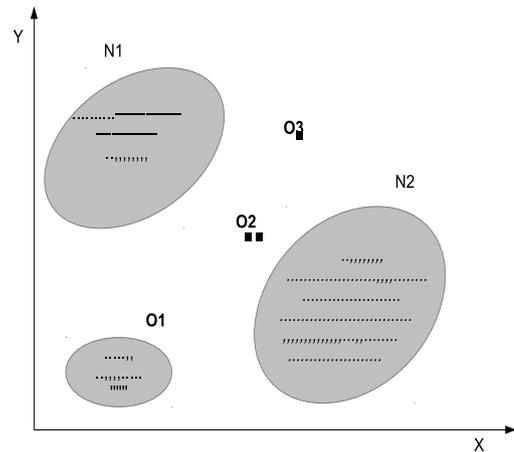


Figure 8: Outlier Detection

mining algorithms involving in two phases:

1. The Training Phase
2. The Detection Phase

During training phase, profiles are created by grouping normal access behaviors and are forwarded in a Batch mode to the Feature Extractor, Feature selector and Classifier. The Classifier produces a trained model out of normal access behavior[108]. Every new test sample during the Detection phase, is made to go

through the same modules: Feature Extractor and Feature Selector, that is finally evaluated by the already trained Classifier. When the sample is found to be deviating from normal profiles, an alarm is raised. The profiles are required to be updated at regular intervals of time and Classifier training is also carried out periodically, so as to minimize the false alarm rate. For Feature selection, we can either employ the Ranking methods or the Filter methods. The Ranking methods output the feature set sorted in descending order according to a particular evaluation measure. The top variables in the feature set are considered to be the most discriminant features. It is therefore essential to determine a threshold to discard features that are considered to have little or no contribution to the classification process. Information Gain(IG) is one of the commonly used evaluation measures.

A variant of IG, with improvisation is the Gain Ratio (GR).

The GR overcomes the bias found in IG towards features resulting in a smaller set of features. For the purpose of Feature Selection we can employ a ranking method that is unsupervised called Principal components analysis(PCA).

The advantage of Filter methods for Feature Selection is that they automatically choose a set of selected features based on a particular evaluation measure. One of the widely employed Filtering methods for Feature Selection is the Best First Search(BFS). It makes use of Forward Selection and Backward Elimination to search through the feature space adopting a Greedy approach. When performance is found to be dropping, it backtracks to the previous feature subsets that have better performance and start all over again from there. BFS is computationally expensive for larger sets. Genetic Algorithms[109] is another type of Filtering technique that is considered to be very effective in practice[110].

## IX. MITIGATING CODE INJECTION ATTACKS

A code injection attack typically involves writing of new machine code into the vulnerable programs memory[111], and after exploiting a bug in the program the control is redirected to the new code[112]. The protection technique[113], W+X mitigates this attack by allowing only either a Write or Execute operations on memory but never allows both[114].

The research progress made so far in this regard is summarized in Table 4.

### a) Types of Code Injection

Some of the flavours of Code Injection attacks are: SQL Injection[121], HTML Script Injection[122], Object Injection[123], Remote File Injection[124] and Code Reuse Attacks(CRAs)[125].

### i. SQL Injection

A technique that uses SQL syntax to input commands that can alter read or modify a database is called SQL Injection. Consider for example a web page having a field on it to allow users to enter a password for authentication. The code behind the page usually a script code, will generate a SQL query to verify the matching password entered against the list of user names:

```
SELECT UsrList.Username FROM UsrList
WHERE UsrList.Password = 'Password'
```

The access is granted when the password entered by the user matches the password specified in the query. If the malicious user can inject some valid code ('password' OR '1'='1') in the Password field. An attacker by leaving the password field empty makes the condition "'1'='1'" to become true and gains access to the database.

### ii. HTML Script Injection

An attacker injects malicious code by making use of the <script>and </script>tags, within which he would change the location property of the document by setting it to an injected script.

### iii. Object Injection

PHP allows serialization and deserialization of objects. If an untrustworthy input is allowed into the deserialization function, it is possible to modify existing classes in the program and execute malicious attacks.

### iv. Remote File Injection

Attackers might provide a Remote Infected file name as the path by modifying the path command of the script file to cause the intended destruction[126].

Table 4: Research Progress in Code Injection Attacks

Authors	Algorithm	Performance	Future enhancement
M Graziano et al.(2016)[115]	Emulation-based framework for ROP	Total analysis time of 4 hours with 16-Core Intel E5-2630 (2.3GHz) and 24GB RAM	Reduction in total analysis time.
Mitropoulos et al.(2016)[116]	Contextual Fingerprinting	Overhead of 11.1% on execution time.	Overhead can be reduced.
A Follner et al.(2015)[117]	Dynamic Binary Instrumentation	2.4x overhead, comparable to similar approaches but no false alarms	The overhead can be reduced
G Parmar et al.(2015)[118]	Input based approach		More techniques/tools for SQLi prevention can be explored or created.
L Deng et al.(2015)[119]	Exception Oriented Programming (EOP)	Detection rate of about 90percent	Can be extended to Mac and Windows kernels.
S Gupta et al. (2015)[120]	Cross-Site Scripting Secure Web Application Framework	Ranging from 1.25% to 5.75% based on the type of JSP program	Discovering the techniques of dropping the HTTP response delay and other rule checks of XSS-SAFE without disturbing its efficiency of XSS attack recognition.

v. Code Reuse Attacks

Attacks in which an attacker directs control flow through an already existing code with an erroneous result are called Code Reuse Attacks[127].

Attackers therefore have come out with code-reuse attacks[128], in which a defect in the software is exploited to create a control flow through existing code-base to a malicious end[129]. The Return Into Lib C(RILC)is a type of code-reuse attack [130] where the stack is compromised and the control is transferred to the beginning of an existing library function such as *mprotect()* to create a memory region[131]that allows both write and execution operations on it to bypass  $W+X$ [132]. Such attacks can be efficiently overcome using Data Mining techniques[133]. The source code is checked to find any such flaws and if so the instructions are classified as malicious[134]. Some of the classification Algorithms that can be used in this Regard are Bayesian[135], SVM[136] and Decision Tree[137].

vi. Return Oriented Programming

ROP attacks start when an attacker gains stack control[138] and redirects the control to a small snippet of code called gadget typically ending with a RET instruction[139]. Because attackers gain control over the return addresses[140], they can assign the RET of one gadget to the start of another gadget[141], achieving the desired functionality out of a large finite set of such small gadgets[142]. ROP Attacks inject no code and yet can induce arbitrary behavior in the targeted system [143]. A compiler-based approach has been suggested in [144] to combat any form of ROP. In [145], the authors present in-place code randomization that can be applied directly on third-party software, to mitigate ROP attacks. Buchanan et al., [146], have demonstrated that return-oriented exploits are practical to write, as the complexity of gadget combination is abstracted behind a programming language and compiler. Davi et al.[147] proposed runtime integrity monitoring techniques that use tracking instrumentation of program binaries based on taint analysis and dynamic tracing. In[148] a tool

DROP, that detects ROP malicious code dynamically, is presented.

vii. Jump Oriented Programming

In Jump Oriented Programming(JOP), an attacker links the gadgets using a finite set of indirect JMP instructions[149], instead of RET instructions. A special gadget called a *dispatcher* is used for flow control management among the gadgets[150].

X. CONCLUSION

The purpose of this survey is to explore the importance of Data Mining techniques in achieving security. The paper is limited to few applications such as Privacy Preserving Data Mining (PPDM), Intrusion Detection System(IDS), Phishing Website Classification, Anomaly/Outlier Detection and Mitigation of Code Injection and Reuse Attacks. Some of the Classification and Clustering algorithms are discussed here considering their significance in Intrusion/Anomaly/ Outlier

Detection Techniques. Other basic Data mining techniques such as Feature Extraction, Association Rule Mining and Decision Trees are also discussed, since many researchers have extensively used these techniques for IDS. The Survey could be made more exhaustive by exploring other security applications of Data Mining such as Malware Detection, Spam Detection, Web Mining and Crime Profiling.

REFERENCES RÉFÉRENCES REFERENCIAS

1. D. R. Stinson, "Cryptography: Theory and Practice 3rd Edition,"*Text Book*, 2006.
2. C.-H. Yeh, G. Lee, and C.-Y. Lin, "Robust Laser Speckle Authentication System through Data Mining Techniques," *IEEE Transactions on Industrial Informatics*, vol. 11, no. 2, pp. 505–512, 2015.
3. S. Khan, A. Sharma, A. S. Zamani, and A. Akhtar, "Data Mining for Security Purpose & its Solitude Suggestions," *International Journal of Technology Enhancements and Emerging Engineering Research*, vol. 1, no. 7, pp. 1–4, 2012.

4. Venugopal K R, K G Srinivasa and L M Patnaik, "Soft Computing for Data Mining Applications," *Springer*, 2009.
5. Vasanthakumar G U, Bagul Prajakta , P Deepa Shenoy, Venugopal K R and L M Patnaik, "PIB: Profiling Influential Blogger in Online Social Networks, A Knowledge Driven Data Mining Approach," *11<sup>th</sup> International Multi-Conference on Information Processing (IMCIP)*, vol. 54, pp. 362–370, 2015.
6. H. Zang and J. Bolot, "Anonymization of Location Data Does Not Work: A Large-Scale Measurement Study," *In Proceedings of the 17th Annual International Conference on Mobile Computing and Networking*, pp. 145–156, 2011.
7. R. J. Bayardo and R. Agrawal, "Data Privacy Through Optimal  $k$ -Anonymization," *21st International Conference on Data Engineering (ICDE'05)*, pp. 217–228, 2005.
8. A. Friedman, R. Wolff, and A. Schuster, "Providing  $k$ -Anonymity in Data Mining," *The VLDB Journal*, vol. 17, no. 4, pp. 789–804, 2008.
9. R. Lu, X. Liang, X. Li, X. Lin, and X. Shen, "EPPA: An Efficient and Privacy-Preserving Aggregation Scheme for Secure Smart Grid Communications," *IEEE Transactions on Parallel and Distributed Systems*, vol. 23, no. 9, pp. 1621–1631, 2012.
10. C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating Noise to Sensitivity in Private Data Analysis," *Theory of Cryptography Conference*, pp. 265–284, 2006.
11. M. Siddiqui, M. C. Wang, and J. Lee, "Detecting Internet Worms Using Data Mining Techniques," *Journal of Systemics, Cybernetics and Informatics*, vol. 6, no. 6, pp. 48–53, 2009.
12. Wu, V. Kumar, J. R. Quinlan, J. Ghosh, Q. Yang, H. Motoda, G. J. McLachlan, A. Ng, B. Liu, S. Y. Philip *et al.*, "Top 10 Algorithms in Data Mining," *Knowledge and Information Systems*, vol. 14, no. 1, pp. 1–37, 2008.
13. M. S. Abadeh, J. Habibi, and C. Lucas, "Intrusion Detection using a Fuzzy Genetics-Based Learning Algorithm," *Journal of Network and Computer Applications*, vol. 30, no. 1, pp. 414–428, 2007.
14. K. S. Desale and R. Ade, "Genetic Algorithm Based Feature Selection Approach for Effective Intrusion Detection System," *International Conference on Computer Communication and Informatics (ICCCI)*, pp. 1–6, 2015.
15. WU Bin, LU Tianliang, ZHENG Kangfeng, ZHENG Dongmei and LIN Xing, "Smartphone Malware Detection Model Based on Artificial Immune System," *China Communications*, vol. 11, no. 13, pp. 86–92, 2014.
16. P Deepa Shenoy, Srinivasa K G, Venugopal K R and L M Patnaik, "Dynamic Association Rule Mining using Genetic Algorithms," *Intelligent Data Analysis*, vol. 9, no. 5, pp. 439–453, 2005.
17. P Deepa Shenoy, Srinivasa K G, Venugopal K R and L M Patnaik, "Evolutionary Approach for Mining Association Rules on Dynamic Databases," *7th Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD)2003, Seoul, South Korea*, pp. 325–336, 2003.
18. S. J. Rizvi and J. R. Haritsa, "Maintaining Data Privacy in Association Rule Mining," *Proceedings of the 28th International Conference on Very Large Data Bases*, pp. 682–693, 2002.
19. S. M. Darwish, M. M. Madbouly, and M. A. El-Hakeem, "A Database Sanitizing Algorithm for Hiding Sensitive Multi-level Association Rule mining," *International Journal of Computer and Communication Engineering*, vol. 3, no. 4, pp. 285–293, 2014.
20. J. Vaidya and C. Clifton, "Privacy Preserving Association Rule Mining in Vertically Partitioned Data," *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 639–644, 2002.
21. J. Vaidya and C. Clifton, "Secure Set Intersection Cardinality with Application to Association Rule Mining," *Journal of Computer Security*, vol. 13, no. 4, pp. 593–622, 2005.
22. M. R. B. Diwate and A. Sahu, "Efficient Data Mining in SAMS through Association Rule," *International Journal of Electronics Communication and Computer Engineering*, vol. 5, no. 3, pp. 593–597, 2014.
23. F. Thabtah, P. Cowling, and Y. Peng, "MCAR: Multi-class Classification based on Association Rule," *The 3rd ACS/IEEE International Conference on Computer Systems and Applications*, pp. 33–39, 2005.
24. K. Hu, Y. Lu, L. Zhou, and C. Shi, "Integrating Classification and Association Rule Mining: A Concept Lattice Framework," *International Workshop on Rough Sets, Fuzzy Sets, Data Mining, and Granular-Soft Computing*, pp. 443–447, 1999.
25. M. Hussein, A. El-Sisi, and N. Ismail, "Fast Cryptographic Privacy Preserving Association Rules Mining on Distributed Homogenous Database," *International Conference on Knowledge-Based and Intelligent Information and Engineering Systems*, pp. 607–616, 2008.
26. J. Zhan, S. Matwin, and L. Chang, "Privacy-Preserving Collaborative Association Rule Mining," *IFIP Annual Conference on Data and Applications Security and Privacy*, pp. 153–165, 2005.
27. F. Giannotti, L. V. Lakshmanan, A. Monreale, D. Pedreschi, and H. Wang, "Privacy-Preserving Mining of Association Rules from Outsourced Transaction Databases," *IEEE Systems Journal*, vol. 7, no. 3, pp. 385–395, 2013.
28. K. Ilgun, R. A. Kemmerer, and P. A. Porras, "State Transition Analysis: A Rule-Based Intrusion Detection Approach," *IEEE Transactions on Software Engineering*, vol. 21, no. 3, pp. 181–199, 1995.

30. V. Kumar, H. Chauhan, and D. Panwar, "K-Means Clustering Approach to Analyze NSL-KDD Intrusion Detection Dataset," *International Journal of Soft Computing and Engineering (IJSCE)*, vol. 3, no. 4, pp. 1–4, 2013.
31. M. Taylor, "Data Mining with Semantic Features Represented as Vectors of Semantic Clusters," *Springer-Verlag*, pp. 1–16, 2012.
32. S. S. Shapiro, "Situating Anonymization within a Privacy Risk Model," *2012 IEEE International Systems Conference(SysCon)*, pp. 1–6, 2012.
33. A.C. Yao, "Protocols for Secure Computations," *23rd Annual Symposium on Foundations of Computer Science, 1982. SFCS'08*, pp. 160–164, 1982.
34. A. Ben-David, N. Nisan, and B. Pinkas, "FairplayMP: A System for Secure Multi-Party Computation," *Proceedings of the 15th ACM Conference on Computer and Communications Security*, pp. 257–66, 2008.
35. P. Mohan, A. Thakurta, E. Shi, D. Song, and D. Culler, "GUPT: Privacy Preserving Data Analysis made Easy," *In Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, pp. 349–360, 2012.
36. A. Kiayias, S. Xu, and M. Yung, "Privacy preserving Data Mining within Anonymous Credential Systems," *International Conference on Security and Cryptography for Networks*, pp. 57–76, 2008.
37. L. A. Dunning and R. Kresman, "Privacy Preserving Data Sharing with Anonymous ID Assignment," *IEEE Transactions on Information Forensics and Security*, vol. 8, no. 2, pp. 402–413, 2013.
38. M. Ouda, S. Salem, I. Ali, and E.-S. Saad, "Privacy-Preserving Data Mining in Homogeneous Collaborative Clustering," *International Arab Journal of Information Technology (IAJIT)*, vol. 12, no. 6, pp. 604–612, 2015.
39. J. Vaidya and C. Clifton, "Privacy-Preserving Data Mining: Why, How, and When," *IEEE Security & Privacy*, vol. 2, no. 6, pp. 19–27, 2004.
40. B. Pinkas, "Cryptographic Techniques for Privacy-Preserving Data Mining," *ACM SIGKDD Explorations Newsletter*, vol. 4, no. 2, pp. 12–19, 2002.
41. M. Roughan and Y. Zhang, "Privacy-Preserving Performance Measurements," *In Proceedings of the 2006 SIGCOMM Workshop on Mining Network Data*, pp. 329–334, 2006.
42. W. Du and Z. Zhan, "Using Randomized Response Techniques for Privacy-Preserving Data Mining," *In Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 505–510, 2003.
43. M. Kantarcioglu, C. Clifton et al., "Privacy-Preserving Distributed Mining of Association Rules on Horizontally Partitioned Data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 9, pp. 1026–1037, 2004.
44. J. Zhan, "Privacy-Preserving Collaborative Data Mining," *IEEE Computational Intelligence Magazine*, vol. 3, no. 2, pp. 31–41, 2008.
45. D. Frey, R. Guerraoui, A. Kermarrec, A. Rault, Taïani, Francois and J. Wang, "Hide and Share: Landmark-Based Similarity for Private KNN Computation," *In Proceedings of the 45th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, pp. 263–274, 2015.
46. T. Zhu, P. Xiong, G. Li and W. Zhou, "Correlated Differential Privacy: Hiding Information in Non-IID Data Set," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 2, pp. 229–242, 2015.
47. BK. Samanthula, Y. Elmehdwi and W. Jiang, "K-Nearest Neighbor Classification over Semantically Secure Encrypted Relational Data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 5, pp. 1261–73, 2015.
48. N P Nethravathi, Prashanth G Rao, P Deepa Shenoy, Venugopal K R and Indramma M, "CBTS: Correlation Based Transformation Strategy for Privacy Preserving Data Mining," *2015 IEEE International WIE Conference on Electrical and Computer Engineering (WIECON-ECE), Dhaka, Bangladesh*, pp. 190–194, 2015.
49. N. Mohammed, D. Alhadidi, B. Fung and M. Debbabi, "Secure Two- Party Differentially Private Data Release for Vertically Partitioned Data," *IEEE Transactions on Dependable and Secure Computing*, vol. 11, no. 1, pp. 59–71, 2014.
50. J. Vaidya, B. Shafiq, W. Fan, D. Mehmood and D. Lorenzi, "A Random Decision Tree Framework for Privacy-Preserving Data Mining," *IEEE Transactions on Dependable and Secure Computing*, vol. 11, no. 5, pp. 399–411, 2014.
51. YJ. Lee, "Privacy-preserving Data Mining for Personalized Marketing," *International Journal of Computer Communications and Networks (IJCCN)*, vol. 4, no. 1, pp. 1–7, 2014.
52. N. Zhang, M. Li, and W. Lou, "Distributed Data Mining with Differential Privacy," *IEEE International Conference on Communications*, pp. 1–5, 2011.
53. F. McSherry and I. Mironov, "Differentially Private Recommender Systems: Building Privacy Into the Net," *In Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 627–636, 2009.
54. A. Friedman and A. Schuster, "Data Mining with Differential Privacy," *In Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 493–502, 2010.
55. M. Roughan and Y. Zhang, "Secure Distributed Data Mining and Its Application to Large Scale Network Measurements," *ACM SIGCOMM Computer Communication Review*, vol. 36, no. 1, pp. 7–14, 2006.

56. N P Nethravathi, Vaibhav J Desai, P Deepa Shenoy, M Indiramma and Venugopal K R, "A Brief Survey on Privacy Preserving Data Mining Techniques," *Data Mining and Knowledge Engineering*, vol. 8, no. 9, pp. 267–273, 2016.
57. A. Narayanan and V. Shmatikov, "De-Anonymizing Social Networks," *IEEE Symposium on Security and Privacy*, pp. 173–187, 2009.
58. S. Chourse and V. Richhariya, "Survey Paper on Intrusion Detection Using Data Mining Techniques," *International Journal of Emerging Technology and Advanced Engineering, ISO*, vol. 4, no. 8, pp. 653–657, 2008.
59. S. Kumar and E. H. Spafford, "A Software Architecture to Support Misuse Intrusion Detection," *Computer Science Technical Report, Purdue University*, pp. 1–19, 1995.
60. J. Allen, A. Christie, W. Fithen, J. McHugh, and J. Pickel, "State of the Practice of Intrusion Detection Technologies," *Technical Report*, pp. 1–239, 2000.
61. R. Mitchell and R. Chen, "Adaptive Intrusion Detection of Malicious Unmanned Air Vehicles Using Behavior Rule Specifications," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 44, no. 5, pp. 593–604, 2014.
62. J. Jabez and B. Muthukumar, "Intrusion Detection System: Time Probability Method and Hyperbolic Hopfield Neural Network," *Journal of Theoretical & Applied Information Technology*, vol. 67, no. 1, pp. 65–77, 2014.
63. E. K. P G Reddy, M. laeng, V. Reddy, and Rajulu, "A Study of Intrusion Detection in Data Mining," *World Congress on Engineering (WCE)*, pp. 6–8, 2011.
64. W. Lee, S. J. Stolfo, and K. W. Mok, "A Data Mining Framework for Building Intrusion Detection Models," *Proceedings of the 1999 IEEE Symposium on Security and Privacy*, pp. 120–132, 1999.
65. S. K. Sahu, S. Sarangi, and S. K. Jena, "A Detail Analysis on Intrusion Detection Datasets," *IEEE International on Advance Computing Conference(IACC)*, pp. 1348–1353, 2014.
66. A. A. C´ardenas, R. Berthier, R. B. Bobba, J. H. Huh, J. G. Jetcheva, D. Grochocki, and W. H. Sanders, "A Framework for Evaluating Intrusion Detection Architectures in Advanced Metering Infrastructures," *IEEE Transactions on Smart Grid*, vol. 5, no. 2, pp. 906–915, 2014.
67. Y. Al-Nashif, A. A. Kumar, S. Hariri, Y. Luo, F. Szidarovsky, and G. Qu, "Multi-Level Intrusion Detection System," *International Conference on Autonomic Computing ICAC'08*, pp. 131–140, 2008.
68. D. Beaver, S. Micali, and P. Rogaway, "The Round Complexity of Secure Protocols," *Proceedings of the 22nd Annual ACM Symposium on Theory of Computing*, pp. 503–513, 1990.
69. W. Lu and I. Traore, "Detecting New Forms of Network Intrusion Using Genetic Programming," *Computational Intelligence*, vol. 20, no. 3, pp. 475–494, 2004.
70. T. F. Lunt, A. Tamaru, and F. Gillham, "A Real-Time Intrusion- Detection Expert System (IDES)," *Technical Report*, pp. 1–166, 1992.
71. D. Barbara, J. Couto, S. Jajodia, L. Popyack, and N. Wu, "ADAM: Detecting Intrusions by Data Mining," *Proceedings of the IEEE Workshop on Information Assurance and Security*, pp. 11–16, 2001.
72. M. Shetty and N. Shekoker, "Data Mining Techniques for Real Time Intrusion Detection Systems," *International Journal of Scientific & Engineering Research*, vol. 3, no. 4, pp. 1–7, 2012.
73. W. Lee, S. J. Stolfo, and K. W. Mok, "Adaptive Intrusion Detection: A Data Mining Approach," *Artificial Intelligence Review*, vol. 14, no. 6, pp. 533–567, 2000.
74. E. Bloedorn, A. D. Christiansen, W. Hill, C. Skorupka, L. M. Talbot, and J. Tivel, "Data Mining for Network Intrusion Detection: How to Get Started," *MITRE*, pp. 1–9, 2001.
75. R. Gopalakrishna and E. H. Spafford, "A Framework for Distributed Intrusion Detection using Interest Driven Cooperating Agents," *Technical Report*, pp. 1–24, 2001.
76. B. Mukherjee, L. T. Heberlein, and K. N. Levitt, "Network Intrusion Detection," *IEEE Network*, vol. 8, no. 3, pp. 26–41, 1994.
77. R. Heady, G. F. Luger, A. Maccabe, and M. Servilla, "The Architecture of a Network Level Intrusion Detection System," *Academic Work submitted to the University of New Mexico*.
78. D. Barbara, N. Wu, and S. Jajodia, "Detecting Novel Network Intrusions using Bayes Estimators." *SDM*, pp. 1–17, 2001.
79. M. Roesch et al., "SNORT: Lightweight Intrusion Detection for Networks," *Proceedings of LISA '99: 13th Systems Administration Conference*, pp. 229–238, 1999.
80. R. Mitchell and R. Chen, "Effect of Intrusion Detection and Response on Reliability of Cyber Physical Systems," *IEEE Transactions on Reliability*, vol. 62, no. 1, pp. 199–210, 2013.
81. D.-Y. Yeung and Y. Ding, "Host-Based Intrusion Detection Using Dynamic and Static Behavioral Models," *Pattern Recognition*, vol. 36, no. 1, pp. 229–243, 2003.
82. W. Lee, S. J. Stolfo, and K. W. Mok, "Mining Audit Data to Build Intrusion Detection Models," *KDD-98 Proceedings*, pp. 66–72, 1998.
83. S. Tanachaiwiwat and K. Hwang, "Differential Packet Filtering Against DDoS Flood Attacks," *ACM Conference on Computer and Communications Security (CCS)*, pp. 1–15, 2003.
84. C.-Y. Tseng, P. Balasubramanyam, C. Ko, R. Limprasittiporn, J. Rowe, and K. Levitt, "A Specification-Based Intrusion Detection System

- forAODV," *Proceedings of the 1st ACM Workshop on Security of Ad-hoc and Sensor Networks*, pp. 125–134, 2003.
85. M. S. Vittapu, V. Sunkari, and A. Y. Abate, "The Practical Data Mining Model for Efficient IDS Through Relational Databases," *International Journal of Research in Engineering and Science*, vol. 3, no. 1, pp. 20–30, 2015.
  86. P. Soni and P. Sharma, "An Intrusion Detection System Based on KDD-99 Data Using Data Mining Techniques and Feature Selection," *International Journal of Soft Computing and Engineering (IJSCE)*, vol. 4, pp. 112–118, 2014.
  87. M. Dubiner, Z. Galil, and E. Magen, "Faster Tree Pattern Matching," *Journal of the ACM (JACM)*, vol. 41, no. 2, pp. 205–213, 1994.
  88. C.-F. Tsai, Y.-F. Hsu, C.-Y. Lin, and W.-Y. Lin, "Intrusion Detection by Machine Learning: A Review," *Expert Systems with Applications*, vol. 36, no. 10, pp. 11 994–12 000, 2009.
  89. D. M. Farid, N. Harbi, and M. Z. Rahman, "Combining Naive Bayes and Decision Tree for Adaptive Intrusion Detection," *arXiv preprint arXiv:1005.4496*, 2010.
  90. A. J. Smola and B. Schölkopf, "A Tutorial on Support Vector Regression," *Statistics and Computing*, vol. 14, no. 3, pp. 199–222, 2004.
  91. T. Joachims, "Text Categorization with Support Vector Machines: Learning With Many Relevant Features," *European Conference on Machine Learning*, pp. 137–142, 1998.
  92. X. Xu and X. Wang, "An Adaptive Network Intrusion Detection Method Based on PCA and Support Vector Machines," *International Conference on Advanced Data Mining and Applications*, pp. 696–703, 2005.
  93. Vasanthakumar G U, P Deepa Shenoy, Venugopal K R and L M Patnaik, "PFU: Profiling Forum Users in Online Social Networks, A Knowledge Driven Data Mining Approach," *2015 IEEE International WIE Conference on Electrical and Computer Engineering (WIECONECE)*, pp. 57–60, 2015.
  94. A. Herzberg and A. Gbara, "Trustbar: Protecting (even naive) Web Users from Spoofing and Phishing Attacks," *Cryptology ePrint Archive Report 2004/155*. <http://eprint.iacr.org/2004/155>, 2004.
  95. R. M. Mohammad, F. Thabtah, and L. McCluskey, "Intelligent Rule- Based Phishing Websites Classification," *IET Information Security*, vol. 8, no. 3, pp. 153–160, 2014.
  96. S. Marchal, J. Francois, R. State, and T. Engel, "PhishStorm: Detecting Phishing with Streaming Analytics," *IEEE Transactions on Network and Service Management*, vol. 11, no. 4, pp. 458–471, 2014.
  97. N. Abdelhamid, A. Ayesh, and F. Thabtah, "Phishing Detection Based Associative Classification Data Mining," *Expert Systems with Applications*, vol. 41, no. 13, pp. 5948–5959, 2014.
  98. Y. Zhang, S. Egelman, L. Cranor, and J. Hong, "Phinding Phish: Evaluating Anti-Phishing Tools," *Academic Work Submitted to School of Computer Science at Research Showcase @ CMU*.
  99. M. Aburrous, M. A. Hossain, K. Dahal, and F. Thabtah, "Predicting Phishing Websites using Classification Mining Techniques with experimental Case Studies," *Seventh International Conference on Information Technology: New Generations (ITNG)*, pp. 176–181, 2010.
  100. J. Chen and C. Guo, "Online Detection and Prevention of Phishing Attacks," *First International Conference on Communications and Networking in China*, pp. 1–7, 2006.
  101. A. Y. Fu, L. Wenyin, and X. Deng, "Detecting Phishing Web Pages with Visual Similarity Assessment Based on Earth Mover's Distance (emd)," *IEEE Transactions on Dependable and Secure Computing*, vol. 3, no. 4, pp. 301–311, 2006.
  102. T. Moore and R. Clayton, "An Empirical Analysis of the Current State of Phishing Attack and Defence," *Academic work*, 2007.
  103. J. Hong, "The State of Phishing Attacks," *Communications of the ACM*, vol. 55, no. 1, pp. 74–81, 2012.
  104. Asha S Manek, P Deepa Shenoy, M Chandra Mohan and Venugopal K R, "Detection of Fraudulent and Malicious Websites by Analysing User Reviews for Online Shopping Websites," *International Journal of Knowledge and Web Intelligence*, vol. 5, no. 3, pp. 171–189, 2016.
  105. C. Jackson, D. R. Simon, D. S. Tan, and A. Barth, "An Evaluation of Extended Validation and Picture-in-Picture Phishing attacks," *International Conference on Financial Cryptography and Data Security*, pp. 281–293, 2007.
  106. C. Cortes, X. Gonzalvo, V. Kuznetsov, M. Mohri, S. Yang, "AdaNet: Adaptive Structural Learning of Artificial Neural Networks," *arXiv:1607.01097v1*, vol. 1, no. 17, pp. 1–18, 2016.
  107. D. Bui, T. Tuan, H. Klempe, B. Pradhan, I. Revhaug, "Spatial Prediction Models for Shallow Landslide Hazards: A Comparative Assessment of the Efficacy of Support Vector Machines, Artificial Neural Networks, Kernel Logistic Regression, and Logistic Model Tree," *springer-Verlag Berlin Heidelberg*, vol. 13, no. 2, pp. 361–378, 2015.
  108. M. Qin and K. Hwang, "Effectively Generating Frequent Episode Rules for Anomaly-based Intrusion Detection," *IEEE Symposium on Security and Privacy*, 2003.
  109. Rasheed and R. Alhajj, "A Framework for Periodic Outlier Pattern Detection in Time-Series Sequences," *IEEE Transactions on Cybernetics*, vol. 44, no. 5, pp. 569–582, 2014.

110. K. Hwang, P. Dave, and S. Tanachaiwiwat, "Netshield: Protocol Anomaly Detection with Datamining against DDOS Attacks," *Proceedings of the 6th International Symposium on Recent Advances in Intrusion Detection, Pittsburgh, PA*, pp. 8–10, 2003.
111. X. Liu, P. Zhu, Y. Zhang, and K. Chen, "A Collaborative Intrusion Detection Mechanism Against False Data Injection Attack in Advanced Metering Infrastructure," *IEEE Transactions on Smart Grid*, vol. 6, no. 5, pp. 2435–2443, 2015.
112. M. G. Schultz, E. Eskin, F. Zadok, and S. J. Stolfo, "Data Mining Methods for Detection of New Malicious Executables," *In Proceedings of IEEE Symposium on Security and Privacy S&P*, pp. 38–49, 2001.
113. B. Wu, T. Lu, K. Zheng, D. Zhang, and X. Lin, "Smartphone Malware Detection Model Based on Artificial Immune System," *China Communications*, vol. 11, no. 13, pp. 86–92, 2015.
114. D. K. B. Patel and S. H. Bhatt, "Implementnig Data Mining for Detection of Malware from Code," *An International Journal of Advanced Computer Technology:Compusoft*, vol. 3, no. 4, pp. 732–740, 2014.
115. Mariano Graziano, Davide Balzarotti and Alain Zidouemba, "ROPMEMU: A Framework for the Analysis of Complex Code-Reuse Attacks," *11th ACM Asia Conference on Computer and Communications Security*, pp. 1–12, 2016.
116. D Mitropoulos, K Stroggylos and D Spinellis, "How to Train Your Browser: Preventing XSS Attacks Using Contextual Script Fingerprints," *ACM Transactions on Privacy and Security*, vol. 19, no. 1, pp. 1–31, 2016.
117. Follner, E. Bodden, "ROPocop - Dynamic Mitigation of Code- Reuse Attacks," *Secure Software Engineering Group*, vol. 29, no. 3, pp. 16–26, 2015.
118. G Parmar and Dr. Kirti Mathur, "Proposed Preventive measures and Strategies Against SQL injection Attacks," *Indian Journal of Applied Research*, vol. 5, no. 5, pp. 716–718, 2015.
119. L. Deng, Q. Zeng, "Exception-Oriented Programming: Retrofitting Code-Reuse Attacks to Construct Kernel Malware," *The Institution of Engineering and Technology*, vol. 10, no. 6, pp. 418–424, 2016.
120. S. Gupta B.B. Gupta, "XSS-SAFE:A Server-Side Approach to Detect and Mitigate Cross-Site Scripting (XSS) Attacks in JavaScript Code," *Springer*, vol. 4, no. 3, pp. 897–920, 2015.
121. M. Polychronakis, "Generic Detection of Code Injection Attacks using Network-Level Emulation," *Ph.D. Thesis*, 2009.
122. P. Rauzy and S. Guilley, "A Formal Proof of Countermeasures Against Fault Injection Attacks on CRT-RSA," *Journal of Cryptographic Engineering*, vol. 4, no. 3, pp. 173–185, 2014.
123. S. Bhatkar, D. C. DuVarney, and R. Sekar, "Address Obfuscation: An Efficient Approach to Combat a Broad Range of Memory Error Exploits," *Usenix Security*, vol. 3, pp. 105–120, 2003.
124. E. G. Barrantes, D. H. Ackley, T. S. Palmer, D. Stefanovic, and D. D. Zovi, "Randomized Instruction Set Emulation to Disrupt Binary Code Injection Attacks," *Proceedings of the 10th ACM Conference on Computer and Communications Security*, pp. 281–289, 2003.
125. S. Bhatkar, D. C. DuVarney, and R. Sekar, "Efficient Techniques for Comprehensive Protection from Memory Error Exploits," *Proceedings of the 14th USENIX Security Symposium*, 2005.
126. Venugopal K R and Rajkumar Buyya, "Mastering C++," *Tata McGraw- Hill Education*, 2013.
127. J. Habibi, A. Panicker, A. Gupta, and E. Bertino, "DISARM: Mitigating Buffer Overflow Attacks on Embedded Devices," *International Conference on Network and System Security*, pp. 112–129, 2015.
128. M. Kayaalp, T. Schmitt, J. Nomani, D. Ponomarev, and N. A. Ghazaleh, "Signature-Based Protection from Code Reuse Attacks," *IEEE Transactions on Computers*, vol. 64, no. 2, pp. 533–546, 2015.
129. E. G˘oktas, E. Athanasopoulos, M. Polychronakis, H. Bos, and G. Portokalidis, "Size Does Matter:Why Using Gadget-Chain Length to Prevent Code-Reuse Attacks is Hard," *USENIX Security Symposium*, pp. 417–432, 2014.
130. S. Kevin Z, F. Monrose, D. Fabian, D. Lucas, L. Alexandra, S. Christopher, and R. Ahmad, "Just-In-Time Code Reuse: On the Effectiveness of Fine-Grained Address Space L]ayout Randomization," *2013 IEEE Symposium on Security and Privacy*, pp. 574–588, 2013.
131. V. van der Veen, E. G˘oktas, M. Contag, A. Pawlowski, X. Chen, S. Rawat, H. Bos, T. Holz, E. Athanasopoulos, and C. Giuffrida, "A Tough Call: Mitigating Advanced Code-Reuse Attacks at the Binary Level," *IEEE Symposium on Security and Privacy*, pp. 1–20, 2016.
132. E. R. Jacobson, A. R. Bernat, W. R. Williams, and B. P. Miller, "Detecting Code Reuse Attacks with a Model of Conformant Program Execution," *International Symposium on Engineering Secure Software and Systems*, pp. 1–18, 2014.
133. M. Musuvathi, D. Y. Park, A. Chou, D. R. Engler, and D. L. Dill, "CMC: A Pragmatic Approach to Model Checking Real Code," *ACM SIGOPS Operating Systems Review*, vol. 36, no. 5, pp. 75–88, 2002.
134. N. Mohanappriya and R. Rajagopal, "Prediction and Pan Code Reuse Attack by Code Randomization Mechanism and Data Corruption," *Techniques and Algorithms in Emerging Technologies*, pp. 162–168, 2016.
135. D. M. Stanley, "CERIAS Tech Report 2013-19 Improved Kernel Security through Code Validation,

- Diversification, and Minimization,” *Ph.D. Thesis*, 2013.
136. Y. Zhuang, T. Zheng, and Z. Lin, “Runtime Code Reuse Attacks: A Dynamic Framework Bypassing Fine-Grained Address Space Layout Randomization,” *SEKE*, pp. 609–614, 2014.
  137. G. F. Roglia, L. Martignoni, R. Paleari, and D. Bruschi, “Surgically Returning to Randomized Lib (C),” *Computer Security Applications Conference*, pp. 60–69, 2009.
  138. E. Buchanan, R. Roemer, H. Shacham, and S. Savage, “When Good Instructions Go Bad: Generalizing Return-Oriented Programming to RISC,” *Proceedings of the 15th ACM Conference on Computer and Communications Security*, pp. 27–38, 2008.
  139. A. Gupta, S. Kerr, M. S. Kirkpatrick, and E. Bertino, “MARLIN: A Fine Grained Randomization Approach to Defend Against ROP Attacks,” *International Conference on Network and System Security*, pp. 293–306, 2013.
  140. S. Checkoway, L. Davi, A. Dmitrienko, A.-R. Sadeghi, H. Shacham, and M. Winandy, “Return-Oriented Programming Without Returns,” *Proceedings of the 17th ACM Conference on Computer and Communications Security*, pp. 559–572, 2010.
  141. L. Davi, A.-R. Sadeghi, and M. Winandy, “ROP Defender: A Detection Tool to Defend Against Return-Oriented Programming Attacks,” *Proceedings of the 6th ACM Symposium on Information, Computer and Communications Security*, pp. 40–51, 2011.
  142. L. Davi, A. Dmitrienko, A.-R. Sadeghi, and M. Winandy, “Return- Oriented Programming Without Returns on ARM,” *Technical Report HGI-TR-2010-002*, 2010.
  143. R. Roemer, E. Buchanan, H. Shacham, and S. Savage, “Return- Oriented Programming: Systems, languages, and applications,” *ACM Transactions on Information and System Security (TISSEC)*, vol. 15, no. 1, pp. 1–34, 2012.
  144. K. Onarlioglu, L. Bilge, A. Lanzi, D. Balzarotti, and E. Kirda, “G-Free: Defeating Return-Oriented Programming Through Gadget- Less Binaries,” *Proceedings of the 26th Annual Computer Security Applications Conference*, pp. 49–58, 2010.
  145. V. Pappas, M. Polychronakis, and A. D. Keromytis, “Smashing the Gadgets: Hindering Return-Oriented programming using In-Place Code Randomization,” *2012 IEEE Symposium on Security and Privacy*, pp. 601–615, 2012.
  146. E. Buchanan, R. Roemer, H. Shacham, and S. Savage, “When Good Instructions Go Bad: Generalizing Return-Oriented Programming to RISC,” *Proceedings of the 15th ACM Conference on Computer and Communications Security*, pp. 27–38, 2008.
  147. L. Davi, A.-R. Sadeghi, and M. Winandy, “Dynamic Integrity Measurement and Attestation: Towards Defense Against Return-Oriented Programming Attacks,” *Proceedings of the 2009 ACM Workshop on Scalable Trusted Computing*, pp. 49–54, 2009.
  148. P. Chen, H. Xiao, X. Shen, X. Yin, B. Mao, and L. Xie, “Drop:Detecting Return-Oriented Programming Malicious Code,” *International Conference on Information Systems Security*, pp. 163–177, 2009.
  149. F. Yao, J. Chen, and G. Venkataramani, “JOP Alarm:Detecting Jump Oriented Programming Based Anomalies in Applications,” *International Conference on Computer Design*, pp. 467–470, 2013.
  150. T. Bletsch, X. Jiang, V. W. Freeh, and Z. Liang, “Jump-Oriented Programming: A New Class of Code-Reuse Attack,” *Proceedings ACM Symposium on Information, Computer and Communications Security*, pp. 30–40, 2011.
  151. S. Abadeh, A. Fernandez, A. Bawakid, S. Alshomrani and F. Herrera, “On The Combination of Genetic Fuzzy Systems and Pairwise Learning for Improving Detection Rates on Intrusion Detection Systems,” *Journal of Expert Systems with Applications*, vol. 42, no. 1, pp. 193–202, 2015.